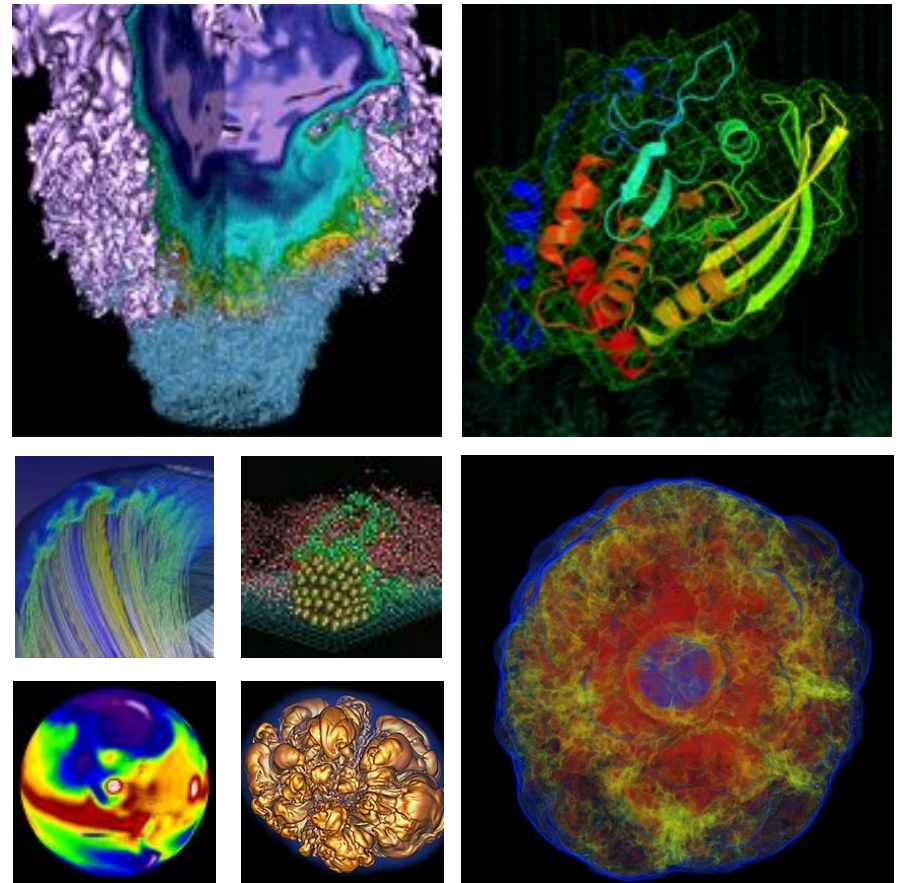


NERSC File Systems Best Practices



Lisa Gerhardt
NERSC Data and Analytics Group

NUG New User Training
February 23, 2017

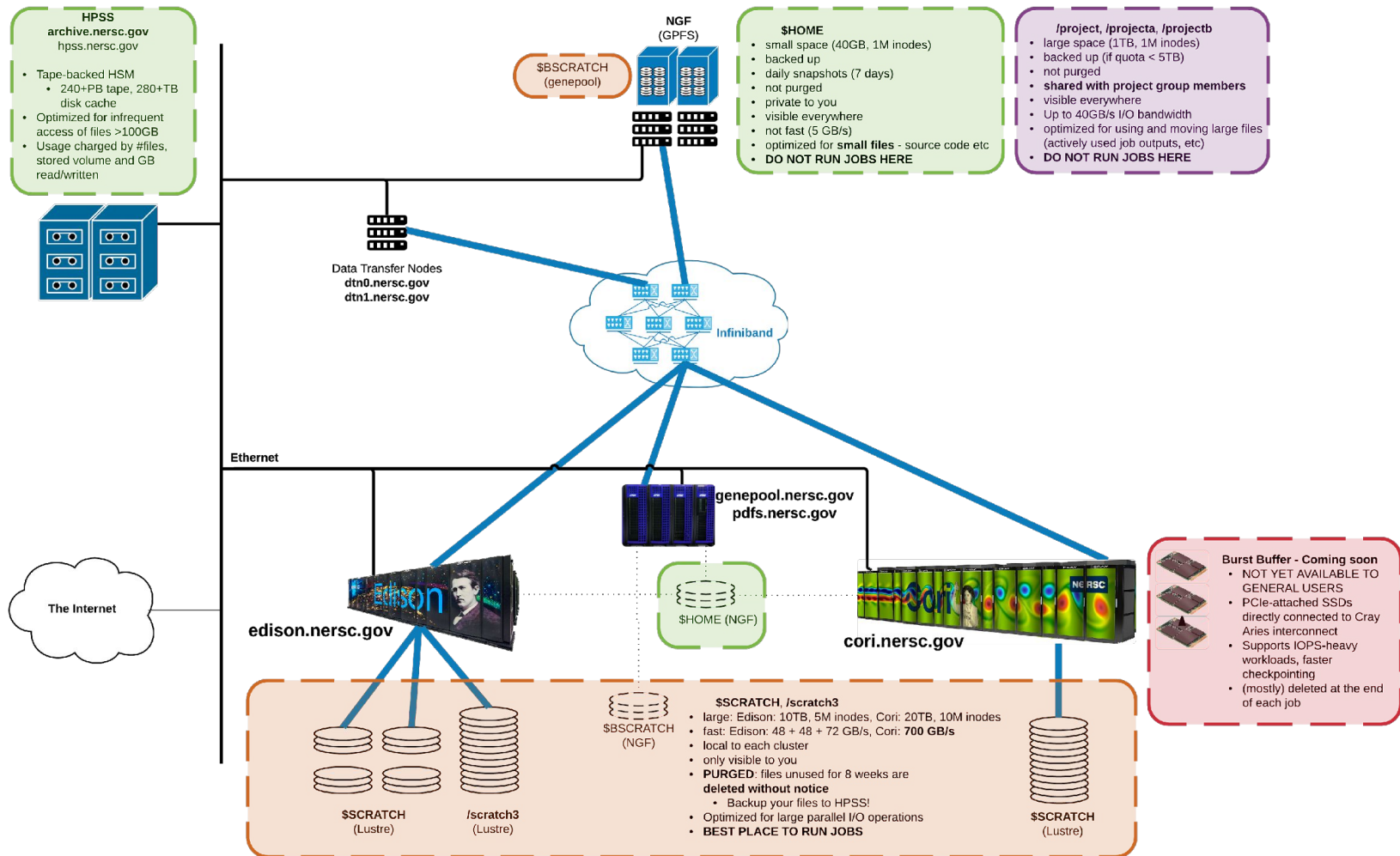
Key Points



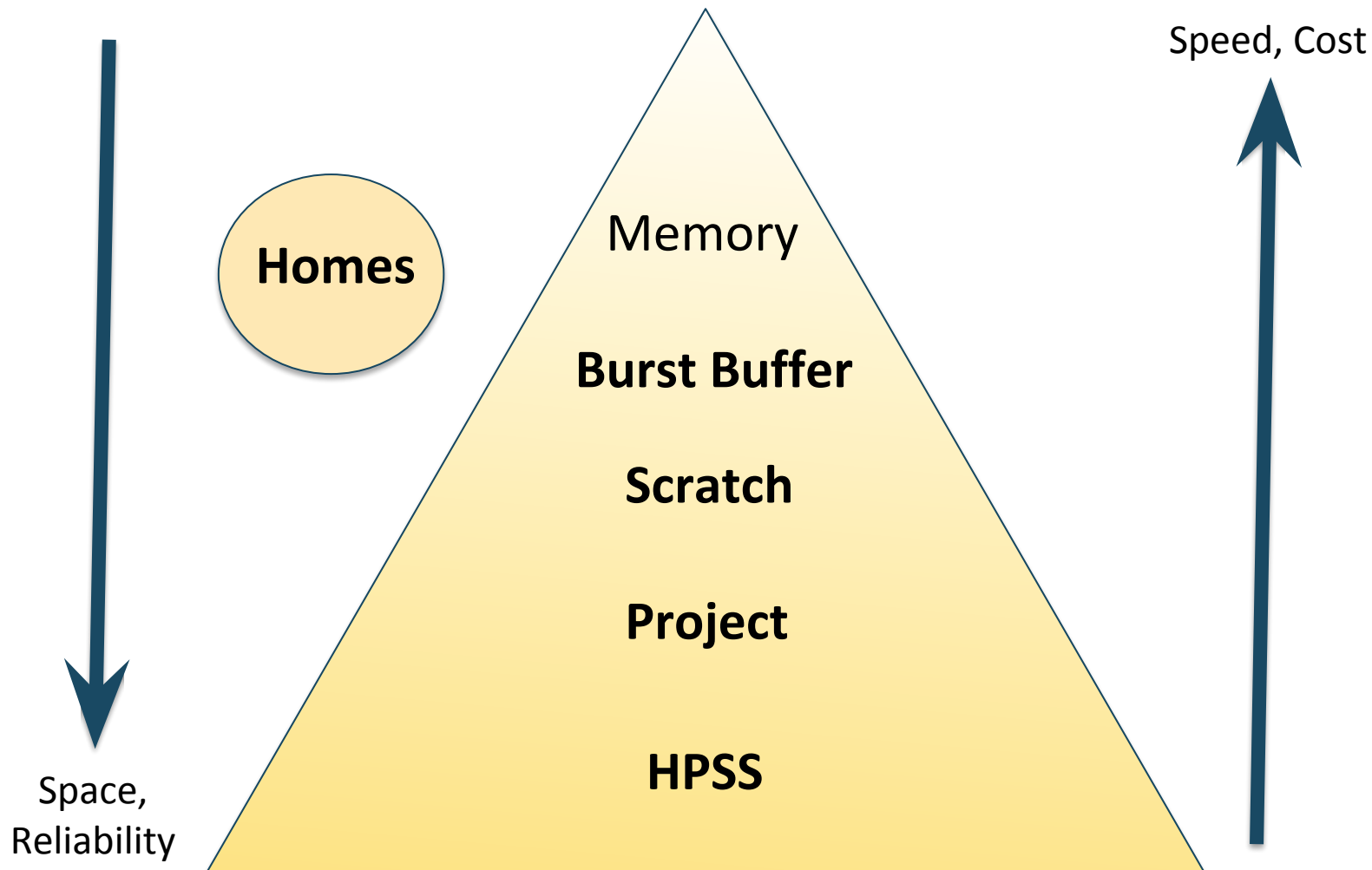
- **Variety of storage types available to meet different needs**
 - Be aware of strengths and limitations of each, use each accordingly
- **BACK UP YOUR IMPORTANT FILES TO HPSS (archive)**
- **If in doubt, ask for help**
 - www.nersc.gov -> “For Users”
 - ServiceNow (help.nersc.gov) or email (consult@nersc.gov)
- **More details tomorrow at the Data Training Day**

NERSC File Systems

NERSC



Simplified NERSC File Systems



- **SSD-equipped nodes (and supporting software) for high-IOPS, high-throughput, “job-local” storage**
 - Directly attached to XC-40 interconnect (Aries)
- **Pre/post-job stage in and stage out**
- **Current configuration:**
 - 288 BB nodes (2 SSDs per BB node)
 - 1.8 PB @ 1.5 TB/s, 12.5M IOPS (measured)
- **USE: Super fast IO layer**
- **More details tomorrow**

Local \$SCRATCH



- **Local to each cluster**
- **Large**
 - Edison: 10 TB, 5M inodes
 - Cori: 20 TB, 1M inodes
 - Monitor with myquota
- **FAST**
 - Edison \$SCRATCH: 168 GB/s aggregate
 - Cori \$SCRATCH: **700** GB/s aggregate
- **Optimized for large parallel I/O workloads**
- **BEST PLACE TO SEND IO FROM JOBS**

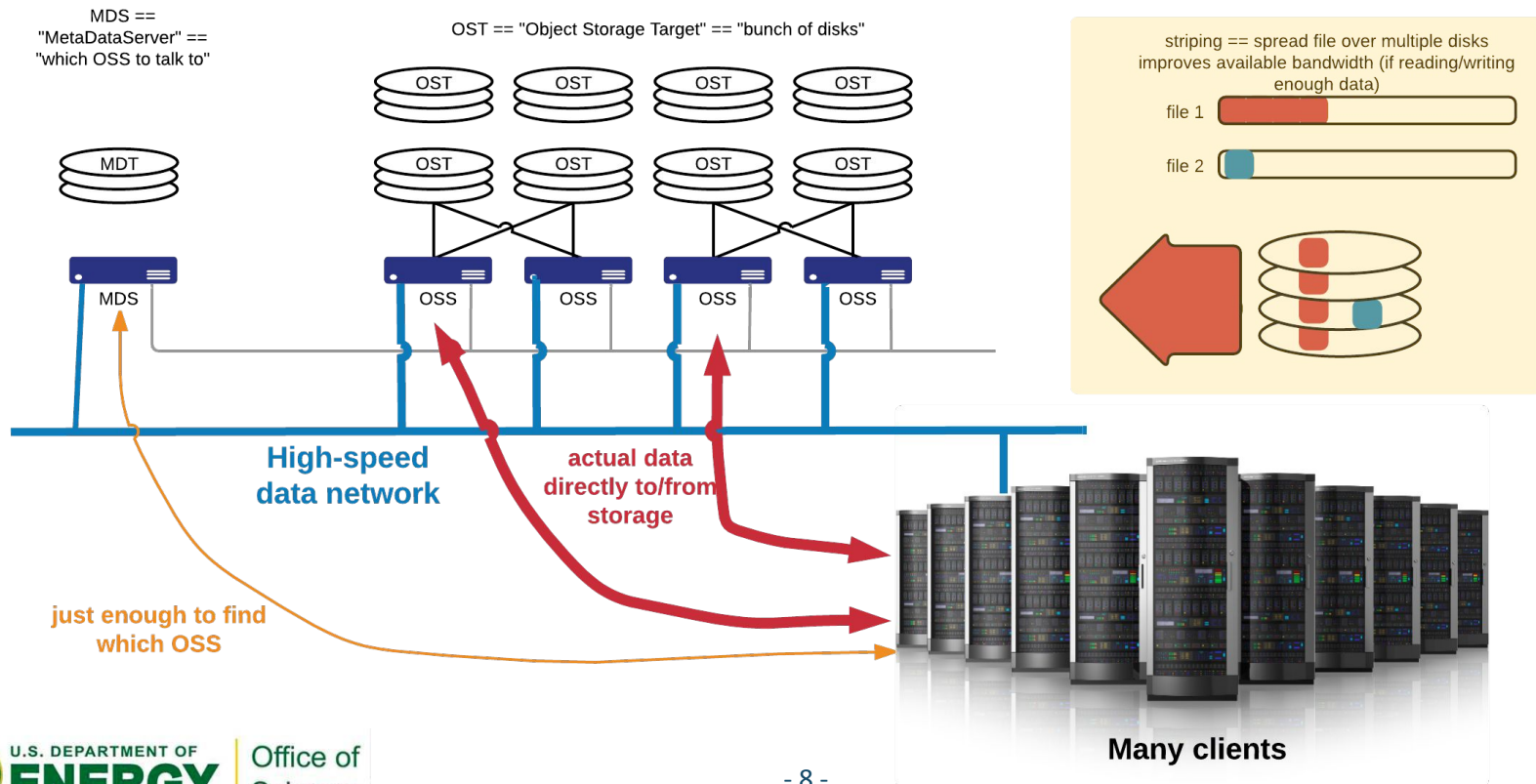
- **Not backed up**
- **Subject to purging**
 - Files not actively used in last 12 weeks (8 weeks on Edison) are **deleted** without notice
 - Purged files are listed in \$SCRATCH/.purged.<timestamp>

BACK UP IMPORTANT FILES TO HPSS!

Local \$SCRATCH



- \$SCRATCH is configured to provide high-bandwidth I/O for many simultaneous users
 - How does it work?



Lustre tips and gotchas



- **Don't keep 100,000 files in the same folder**
 - Hard work for OSS, affects performance for other users
 - 100 folders with 1000 files each is much faster
- **'ls' vs 'ls -l'**
 - Passing options to 'ls' invokes an inquiry on each inode in the folder – occupies OSS/OST with small transfers, non-optimal
 - Basic 'ls' needs only information kept in MDS, much faster
- **'lfs find' vs 'find'**
 - Same principle: special (limited) version of find that only uses data on MDS, not OSS/OST

Project File Systems



- **Large space (1TB, 1M inodes)**
- **USE: holding and sharing actively-used data**
- **Every group gets a project directory**
 - `/project/projectdirs/m9999`
 - Group readable and writeable
- **Daily snapshots for last 7 days**
- **Never purged**
- **Visible everywhere**
 - Web-accessible via *science gateways*
- **Aggregate bandwidth ~150 GB/s, brokered by DVS**
 - Not the place to run jobs .. but jobs could read large input files directly from here

HPSS: Long Term Tape Archive

NERSC



- **archive.nersc.gov**
 - HSM: disk cache, ultimately everything is stored on tape
 - Parallel connections over NERSC internal 10GbE network
- **Available to all NERSC users**
- **No quota, but charged in “Storage Resource Units”**
 - Function of volume of data in storage, number of files in storage and volume of data transferred
 - Monitor usage via NIM
- **USE: Long term storage of data**

Accessing HPSS



| Tool | What it does | Where/why to use it | Example |
|---------------|---------------------------------|--|--|
| htar | Tar directly to/from HPSS | From NERSC hosts. Simple store/retrieve of large directories | <pre>\$ htar cf results-for-publication.tar my_results/</pre> |
| hsi | CLI client | From NERSC hosts. Full featured client | <pre>\$ hsi A:/home/s/sleak-> put myfile</pre> |
| Globus Online | Data transfer service | Fire-and-forget transfers | See www.globusonline.org |
| pftp, ftp | High performance (parallel) ftp | When need/prefer ftp-like interface | <pre>\$ pftp archive.nersc.gov ftp> pput results-for-publication.tar</pre> |
| gridFTP | | External, gridFTP-enabled sites (you need a grid credential) Note: garchive.nersc.gov | <pre>\$ globus-url-copy file://\$HOME/myresults.tar gsiftp://garchive.nersc.gov/home/ s/sleak/results-for-publication.t ar</pre> |

- **Store files as you intend to extract them**
 - Backup to protect against accidental deletion: use htar to bundle up each directory
 - Archiving data mirror: bundle by month data was taken or detector run characteristics, etc.
- **Optimal size of bundles is currently 100s of GB**
 - Larger than >1 TB retrieval is prone to interruption
- **User xfer queue for long running transfer**
 - Archiving during compute job only gets single stream data movement AND costs MPP hours
- **Use two-factor process to transfers in/out of HPSS from outside of NERSC**
 - Globus between centers, hsi/htar to/from HPSS

HPSS is a Tape System



- **All data in HPSS eventually ends up on tape**
 - Transfers in go first to disk cache, so they are very quick
- **Tape is linear media**
 - Data cannot be written anywhere, only appended at end
 - Reading and writing are sequential, not random-access
 - Robot must fetch tape, load it into drive, read forwards until file is reached, then read file
 - Number-of-files has bigger impact on access performance than number-of-GB
- **If you are retrieving more than ~100 files, please order your retrievals by tape position**
 - NERSC has a helper script and instructions to help you sort

Checking my Usage



- **nim.nersc.gov**

My NGF Quotas & Usage

| Username | Full Name | Home Space Used (GiB) | Home Space Quota (GiB) | HSQ Def? | Home Inodes Used | Home Inode Quota | HIQ Def? | Home Quota End | Prop Chng | |
|----------|--------------|-----------------------|------------------------|----------|------------------|------------------|----------|----------------|-----------|------------------------------------|
| sleak | Stephen Leak | 6.1 | 40 | Y | 133,443 | 1,000,000 | Y | Never | N | Update User Quotas |

Usage for My Project Directories

| Project Directory | Owner | Group Name | ERCAP Project | Space Usage | Space Quota | Default Space Quota? | Space% | Inode Usage | Inode Quota | Default Inode Quota? | Inode% | Quota Expiration Date | Projdir Status | Status Effective Date | Projdir ID | Group ID | Project ID | Prop Chng | |
|-------------------|---------|------------|---------------|-------------|-------------|----------------------|--------|-------------|-------------|----------------------|--------|-----------------------|----------------|-----------------------|------------|----------|------------|-----------|-------------------------------------|
| carver | dpaul | mpccc | staff | 8 | 1.0 | Y | 0.8 | 63,918 | 1,000,000 | Y | 6 | Never | Active | Jan-06-2016 | 43906 | 11988 | 13439 | N | View Projdir Quotas |
| dirac | whitney | mpccc | staff | 165 | 1.0 | Y | 16 | 15,576 | 1,000,000 | Y | 1.6 | Never | Active | Jan-06-2016 | 43946 | 11988 | 13439 | N | View Projdir Quotas |
| genepool | jay | mpccc | staff | 130 | 1.0 | Y | 13 | 900,469 | 1,000,000 | Y | 90 | Never | Active | Jan-06-2016 | 43970 | 11988 | 13439 | N | View Projdir Quotas |

- **myquota**
- **prjquota**

- **Home directory shared across all NERSC clusters**
 - Private to you
- **USE: small source code, configuration files, etc**
- **Small space (40GB, 1M inodes)**
 - Check usage with “myquota” command
- **Backed up to tape, and daily snapshots for last 7 days**
- **Never purged**
- **Visible everywhere**
- **Aggregate bandwidth ~50 GB/s**
- **DO NOT RUN JOBS HERE**
 - Don't send Slurm stderr/stdout here

- **\$HOME daily snapshots (last 7 days)**
 - Extra-hidden folder \$HOME/.snapshots

```
sleak@cori03:~$ ls -a
.          .bashrc.ext  .globus    .local     .pyhistory  .udiRoot    .zprofile.ext
..         .cache       .history   .login     .python-eggs .vim         .zshenv
.Xauthority .config      .inputrc   .login.ext .ssh         viminfo     .zshenv.ext
.bash_history .cshrc      .intel     .netrc     .subversion  .vimrc      .zshrc
.bash_profile .cshrc.ext  .kshrc     .odbc.ini  .swp         .zlogin     .zshrc.ext
.bash_profile.ext .fontconfig .kshrc.ext .profile   .tcshrc     .zlogin.ext my_stuff
.bashrc      .gitconfig  .lessht    .profile.ext .tcshrc.ext .zprofile

sleak@cori03:~$ ls .snapshots
2016-03-09  2016-03-10  2016-03-11  2016-03-12  2016-03-13  2016-03-14  2016-03-15  2016-03-16

sleak@cori03:~$ ls .snapshots/2016-03-12
NESAP  Tools  Training  UserSupport  aaa  bin  intel  log.lammps  xtnodestat
```

- **Mistakes, hardware failures happen!**
Backup important files to HPSS

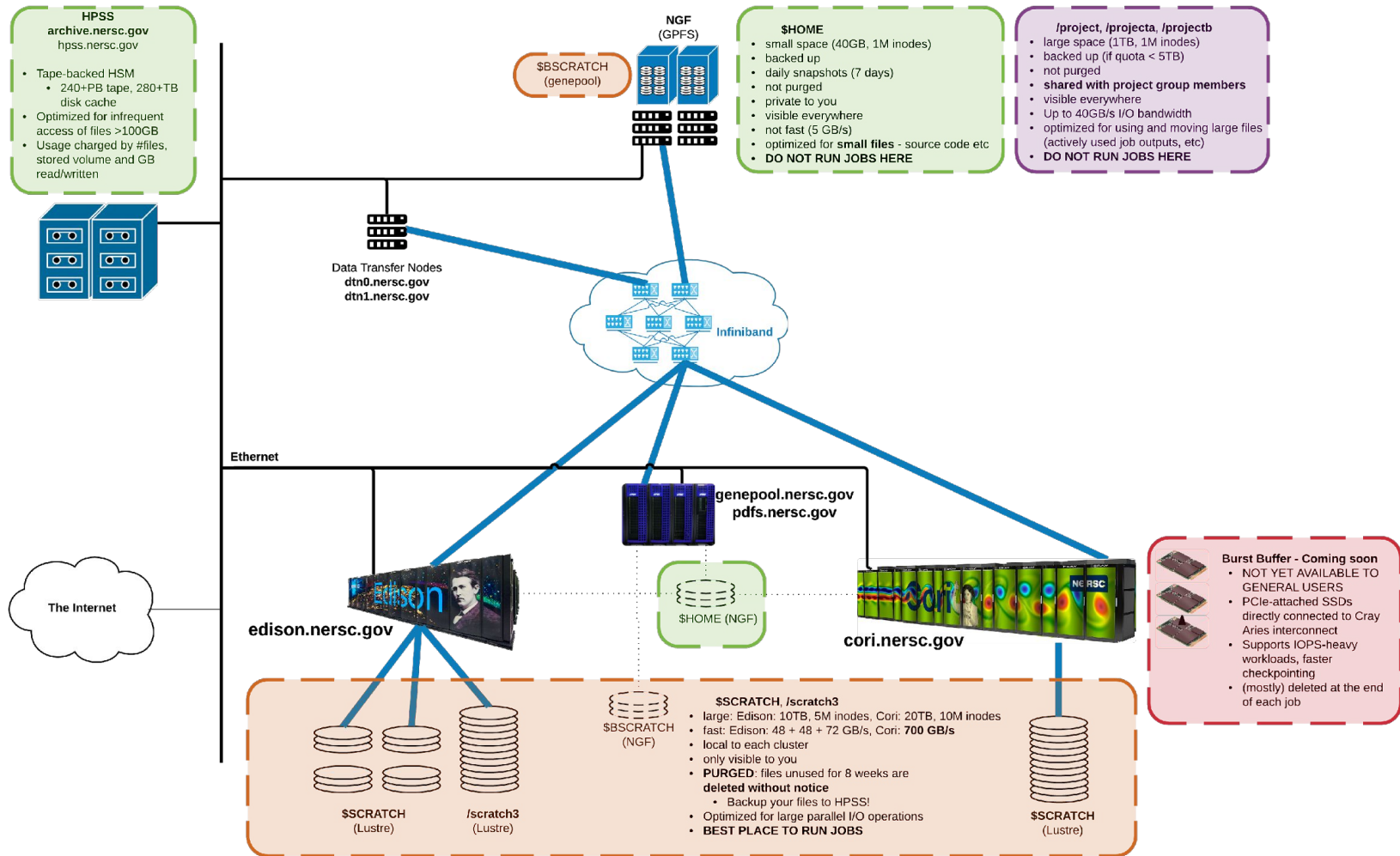
Where do I put my software?



- **Generally, scratch will give best performance**
 - Super fast IO, but is purged
- **If you have a very large software stack with lots of dynamic libraries**
 - encapsulate it with Shifter (details tomorrow)
 - Details on optimizing python based stacks tomorrow

NERSC File Systems Summary

NERSC





National Energy Research Scientific Computing Center

- **Served from NERSC Global Filesystem (NGF)**
 - Based on IBM GPFS
- **Provided by two ~100 TB file systems**
 - /global/u1/
 - /global/u2/
 - Users assigned randomly to one of them
 - Symbolic link on the other
- **Access it with \$HOME or ~/**
 - Underlying name might change, "\$HOME" will not

- **Served from NERSC Global Filesystem (NGF)**
 - Based on IBM GPFS
- **5 GB/s aggregate bandwidth**
 - To \$HOME, shared by all users
- **Shared by ~6000 active NERSC users**
 - Inefficient use affects others
- **Don't run jobs here!**
 - Neither space nor I/O bandwidth are suitable
- **Don't send Slurm stderr/stdout here**
 - Submit jobs from \$SCRATCH, or redirect output to there

- Quotas
 - 40 GB
 - 1,000,000 inodes (i.e. files and directories)
 - Quota increases for \$HOME are almost never granted
 - (why do you need more than 40GB of source code? May need to reconsider what you are storing in \$HOME)
 - Monitor your usage with `myquota`
 - Also visible in NIM

```
sleak@cori03:~$ myquota
```

```
Displaying quota usage for user sleak:
```

| FileSystem | Usage | Space (GB) | Quota | InDoubt | Usage | Inode | Quota | InDoubt |
|------------------|-------|------------|-------|---------|--------|----------|-------|---------|
| /global/cscratch | 0 | | 20480 | - | 51 | 10000000 | | - |
| HOME | 6 | | 40 | 0 | 133431 | 1000000 | | 0 |

- **Help! I deleted some large files, but my usage according to my quota stayed the same**
 - Check for any running processes that are using the deleted files. The space will not be returned until these processes finish or are killed
 - The process may be on a different login node, or part of a batch job you have running

- **Backups and retention**
 - Nightly backups to tape
 - Kept for 90 days
 - Last 7 days accessible via hidden \$HOME/.snapshots folder
 - Recovering from tape is possible but slow, contact us via ServiceNow (help.nersc.gov) or email (consult@nersc.gov)
 - Data is kept on tape for 1 year after your account is deactivated

Project File Systems



- **Served from NERSC Global Filesystem (NGF)**
- **5.1 PB high-performance disk**
 - 50GB/s aggregate bandwidth
- **Every MPP repo has a project space**
 - `/project/projectdirs/m9999`
- **Tuned for large streaming file access**
 - Not the place to run jobs .. But jobs could read large input files directly from here

- **Sharing data**

- Access control is via Unix groups
- PI manages membership
 - (<http://www.nersc.gov/users/accounts/nim/nim-guide-for-pis/>)
- More on sharing soon

- **Science gateways**

- Web portals for sharing data with external collaborators

```
mkdir /project/projectdirs/yourproject/www
```

```
chmod -R 755 /project/projectdirs/yourproject/www
```

- Corresponds to <http://portal.nersc.gov/project/yourproject>
- See <http://www.nersc.gov/users/data-analytics/science-gateways/>

Project File Systems



- Quotas

- 1 TB
- 1,000,000 inodes (i.e. files and directories)
- Quota increases considered
 - <http://www.nersc.gov/users/storage-and-file-systems/file-systems/disk-quota-increase-request/>
- Monitor your usage with `prjquota <yourproject>`
 - Also visible in NIM

```
sleak@cori03:~$ prjquota acme
```

| Project | Usage | Space (GB) Quota | InDoubt | Usage | Inode Quota | InDoubt |
|---------|-------|---------------------|---------|--------|----------------|---------|
| acme | 1014 | 1024 | 0 | 899382 | 1000000 | 0 |

- **Backups and retention**
 - Nightly backups to tape
 - Kept for 90 days
 - Last 7 days accessible via hidden \$HOME/.snapshots folder
 - Recovering from tape is possible but slow, contact us via ServiceNow (help.nersc.gov) or email (consult@nersc.gov)
 - Data is kept on tape for 1 year after project becomes inactive (no allocation, no activity)

- Quotas

- Edison: 10 TB, 5,000,000 inodes
- Cori: 20 TB, 10,000,000 inodes
- Quota increases considered
 - <http://www.nersc.gov/users/storage-and-file-systems/file-systems/disk-quota-increase-request/>
- Monitor your usage with myquota
 - Also visible in NIM

```
sleak@cori03:~$ myquota
```

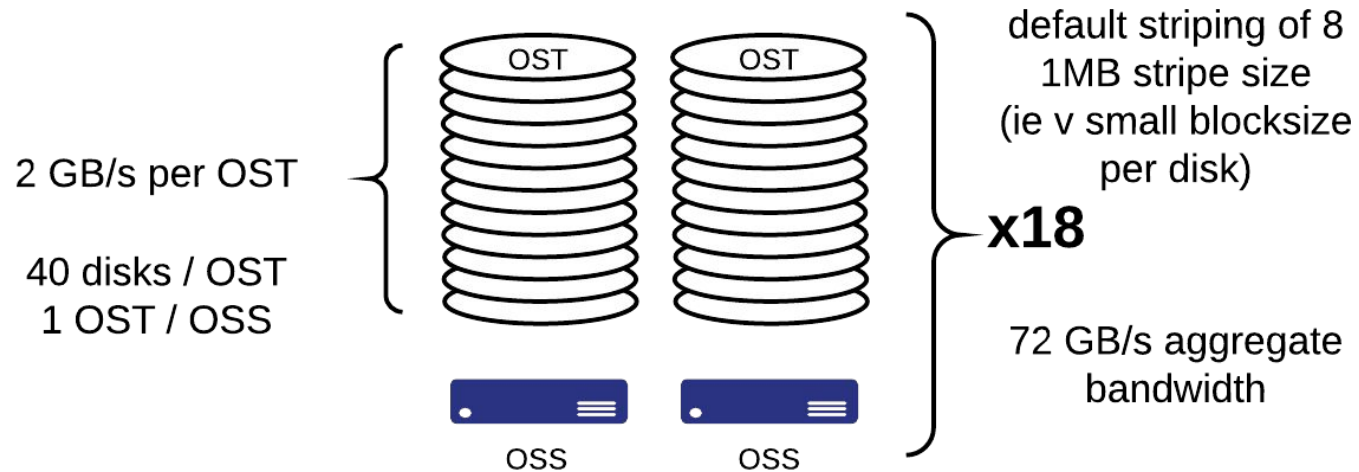
```
Displaying quota usage for user sleak:
```

| FileSystem | Usage | Space (GB) Quota | InDoubt | Usage | Inode Quota | InDoubt |
|------------------|-------|---------------------|---------|--------|----------------|---------|
| /global/cscratch | 0 | 20480 | - | 51 | 10000000 | - |
| HOME | 6 | 40 | 0 | 133431 | 1000000 | 0 |

- **Lustre filesystem**

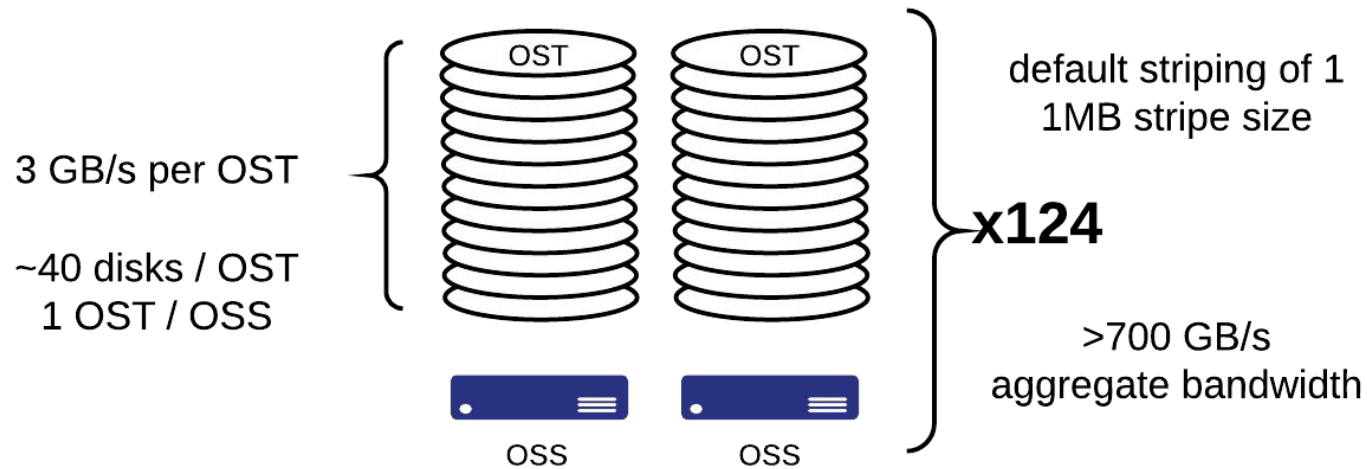
- Edison: provided by two 2 PB filesystems
 - Users assigned randomly to one of them
- Cori: single 28 PB filesystem
- Access it with \$SCRATCH
- Edison /scratch3: access considered by request
 - <http://www.nersc.gov/users/computational-systems/edison/file-storage-and-i-o/>
 - Access it by name (/scratch3/scratchdirs/\$USER)
 - /scratch3 has greater I/O bandwidth





- **I/O striped over 8 OSTs of 40 disks each**
 - high I/O bandwidth

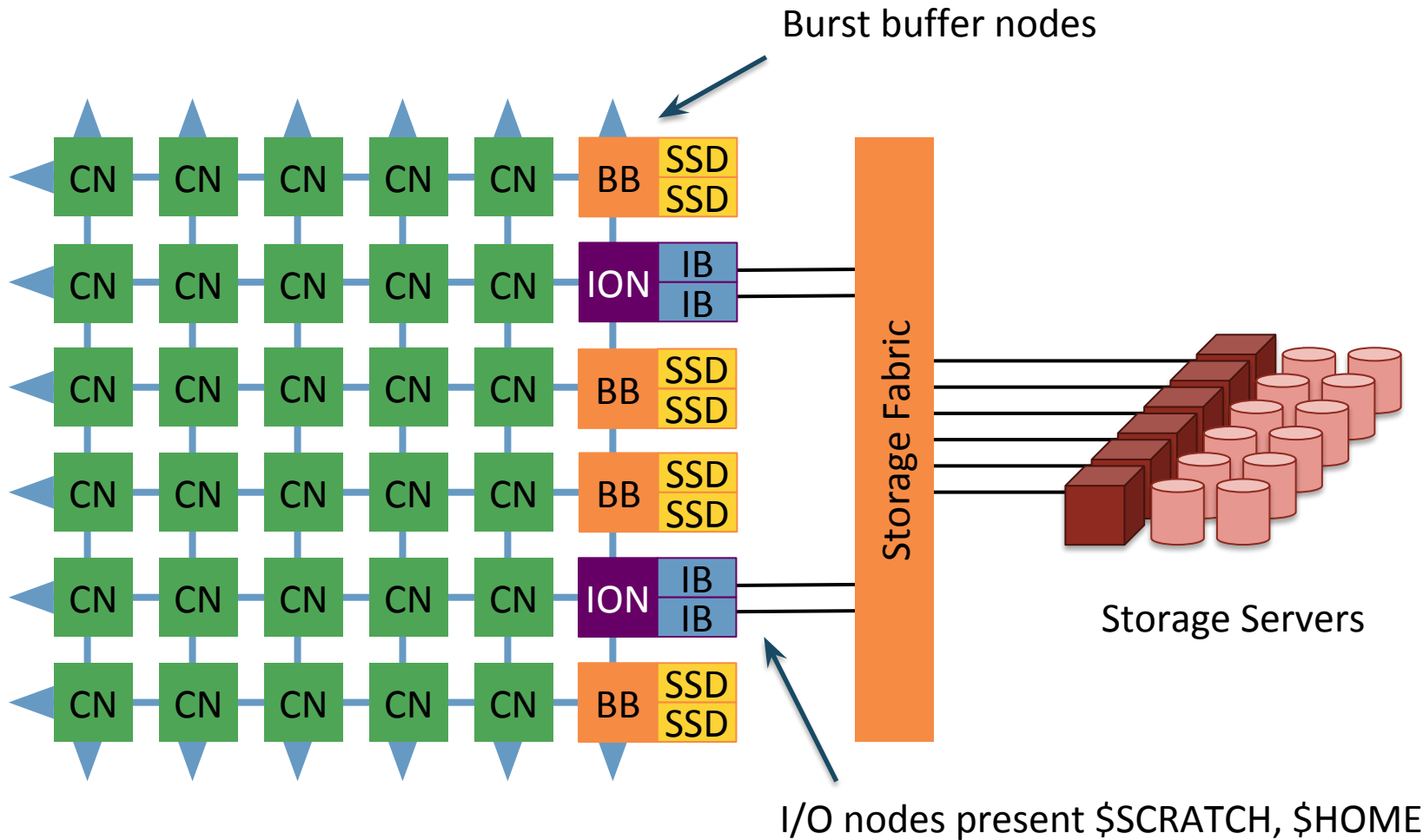
Cori \$SCRATCH



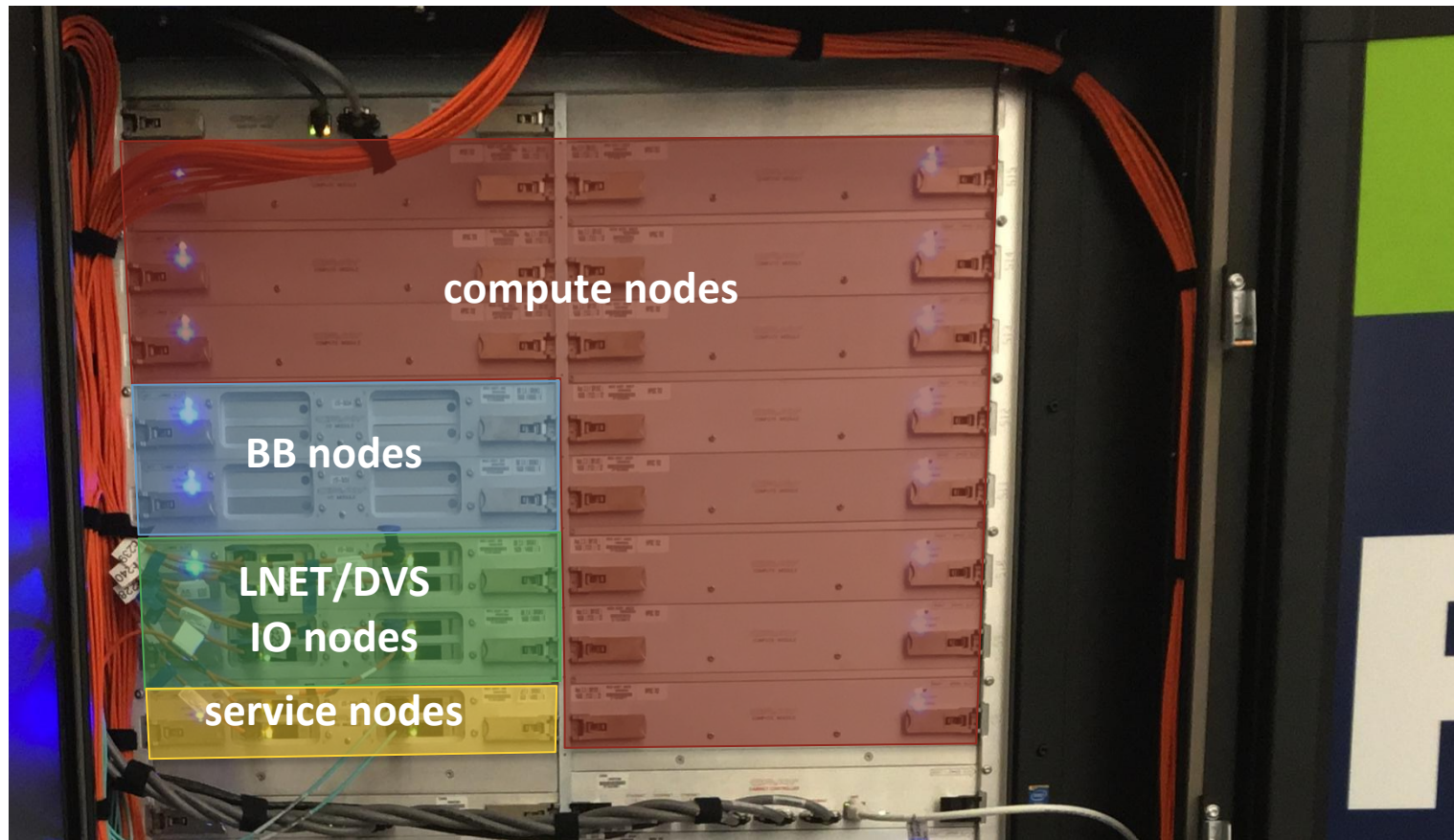
- **Large space, highly parallel**
 - Eventually will become global scratch space

- **Why?**
 - Limitations of \$SCRATCH:
 - Relies on large, throughput-oriented I/O for performance
 - Checkpointing – extreme bandwidth requirements
 - 1000's of nodes each writing 10's of GB
 - Mostly not required again
 - For large parallel jobs, I/O is often “bursty”
 - Most cores waiting while few cores do I/O
- **How?**
 - #BB job directives passed to sbatch

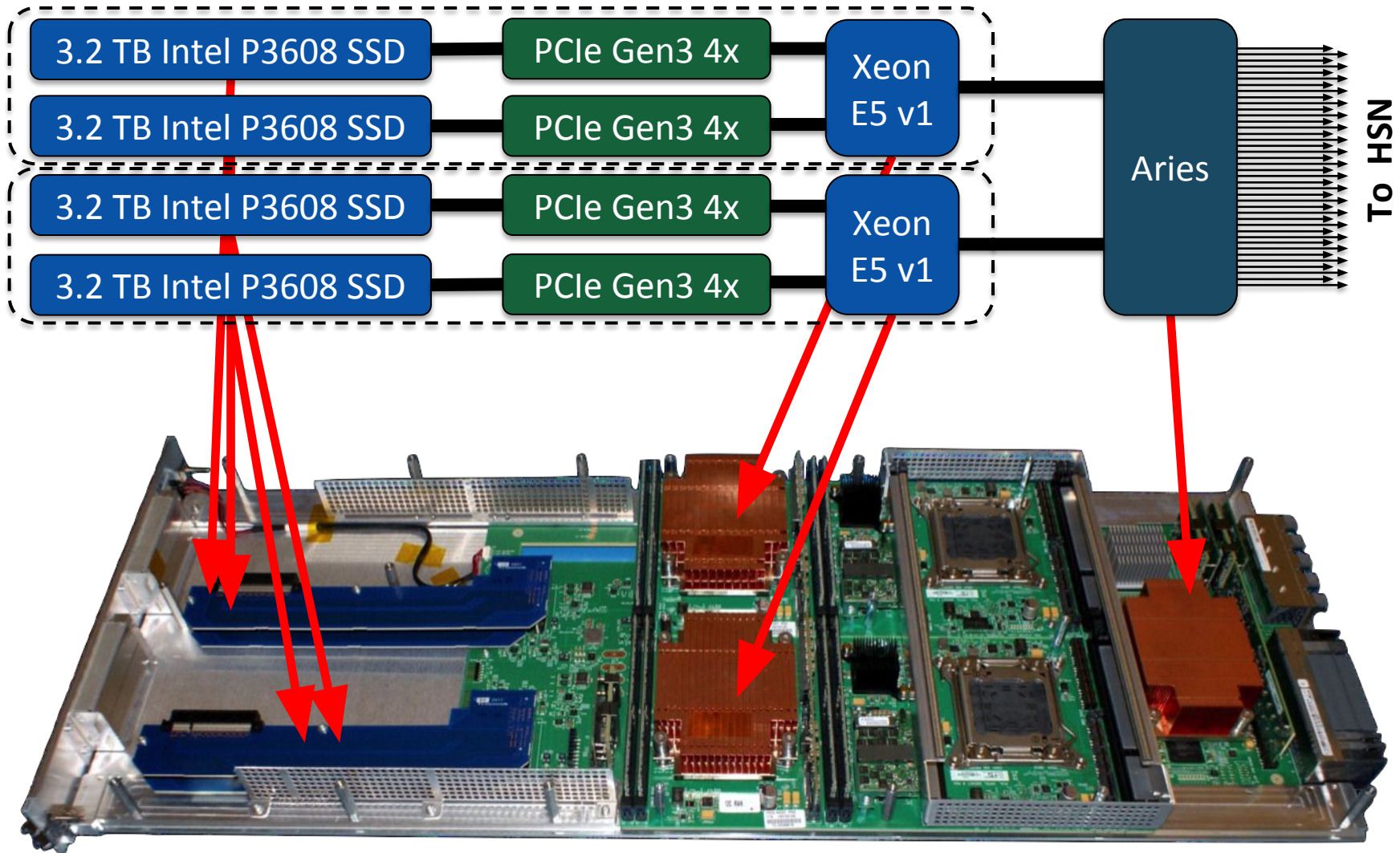
Burst Buffer

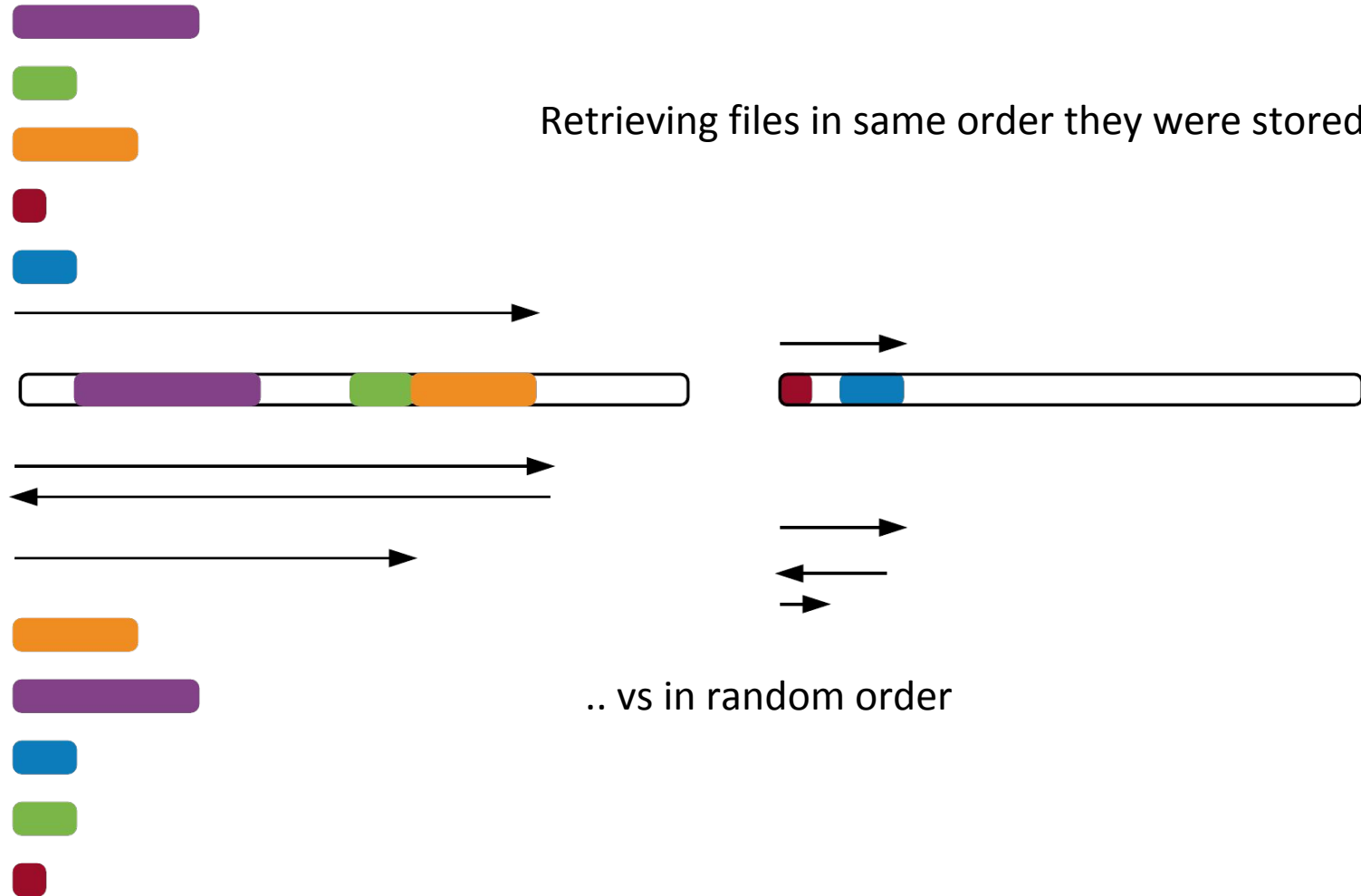


Burst Buffer



Burst Buffer





- **Best practices/Worst practices:**
 - <http://www.nersc.gov/users/storage-and-file-systems/hps/s/storing-and-retrieving-data/mistakes-to-avoid/>
 - Store a few very large files, not many small files
 - htar or tar-first-in-\$SCRATCH
 - Recursively storing or fetching a directory tree will result in many unordered accesses
 - Use htar or tar instead
 - hpss_file_sorter.script => sorts a list of files into “tape order”

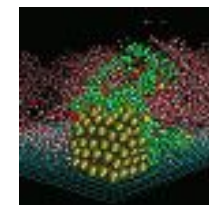
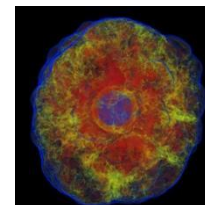
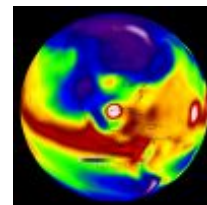
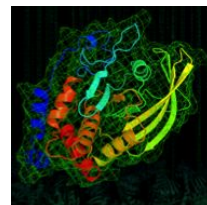
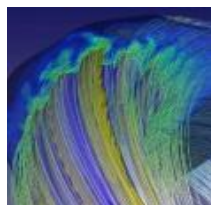
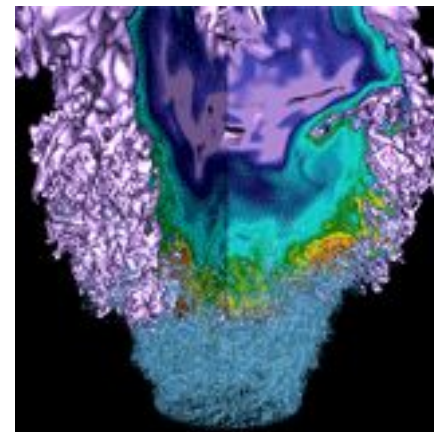
Optimizing I/O Performance



- **You can view/change the stripe size**
 - `lfs getstripe $SCRATCH/my_file.dat`
 - `lfs setstripe -s 4m -c 4 $SCRATCH/my_folder`
 - (affects new files, not existing files)
- **Some shortcuts for single-shared-file I/O:**
 - `stripe_small $SCRATCH/my_folder`
 - Files >1 GB
 - `stripe_medium $SCRATCH/my_folder`
 - Files >10 GB
 - `stripe_large $SCRATCH/my_folder`
 - Files >100 GB
- **Use with care: can make performance worse**
 - <http://www.nersc.gov/users/storage-and-file-systems/optimizing-io-performance-for-lustre/>

- **Best practices/Worst practices:**
 - <http://www.nersc.gov/users/storage-and-file-systems/hps/s/storing-and-retrieving-data/mistakes-to-avoid/>
 - HPSS has a single database instances, all user interactions trigger database activity
 - hsi -q 'ls -l' is database intensive, $O(N^2)$ with number of files in directory
 - Too many files in one folder can lock up system for everybody
 - Streaming data to pftp from Unix pipeline
 - HPSS does not know how big the data will be, likely to put it in wrong place
 - Vulnerable to network glitch

Sharing Data



- **Security matters!**
 - Never share passwords
- **With other NERSC users**
 - Project directories (/project) are designed for sharing files with colleagues
 - Not \$HOME
 - Unix groups, ACLs (“file access control lists”)
 - give, take commands
- **With external collaborators**
 - Science gateways (on /project)

- **Unix groups**

- What groups am I in?
 - `groups`
- New files are associated with your default group
- To change which group the file is associated with:
 - `chgrp my_other_group myfile.txt`
 - `chgrp -R my_other_group whole_directory_tree/`
- To ensure users in `my_other_group` can read/write a file or folder:
 - `chmod g+rw myfile.txt`
 - `chmod g+rws my_new_folder/`
 - “s” – setgid

“setgid” ??



- **setgid “set group id”**
 - File mode, set with `chmod`
 - When set on a folder, it means “things added to this folder should inherit the group of the folder”
 - (so I don’t need to keep typing `chgrp` for each new file)
 - NOTE: only things added, not things that were already there

- **Finer-grain control of access**

- `getfacl`, `setfacl`
- `setfacl -m u_or_g:who:what_perms myfile.txt`
- `setfacl -x`
 - Remove a FACL

```
getfacl some_file.txt
# file: some_file.txt
# owner: sleak
# group: sleak
user::rw-
group::r--
other::---
```

```
$ setfacl -m u:rjhb:rw some_file.txt
$ getfacl some_file.txt
# file: some_file.txt
# owner: sleak
# group: sleak
user::rw-
user:rjhb:rw-
group::r--
mask::rw-
other::---
```

My colleague still can't see my file?



- **Check permissions of the folder it is in, and the folder above that, etc**
 - Missing permissions at any point in the tree will prevent access to the next level of the tree
- **Don't forget "x" on folders**

Give and Take



- **Appropriate for smaller files**

```
joe% give -u bob coolfile
```

- File copied *to* spool location
- Bob gets email telling him Joe has given him a file

```
bob% take -u joe coolfile
```

- File copied *from* spool location

- **Make data available to outside world**

```
mkdir /project/projectdirs/bigsci/www  
chmod o+x /project/projectdirs/bigsci  
chmod o+rx /project/projectdirs/bigsci/www
```

- **Access with web browser**

```
http://portal.nersc.gov/project/bigsci
```

- **More info:**

- <https://www.nersc.gov/users/data-analytics/science-gateways/>

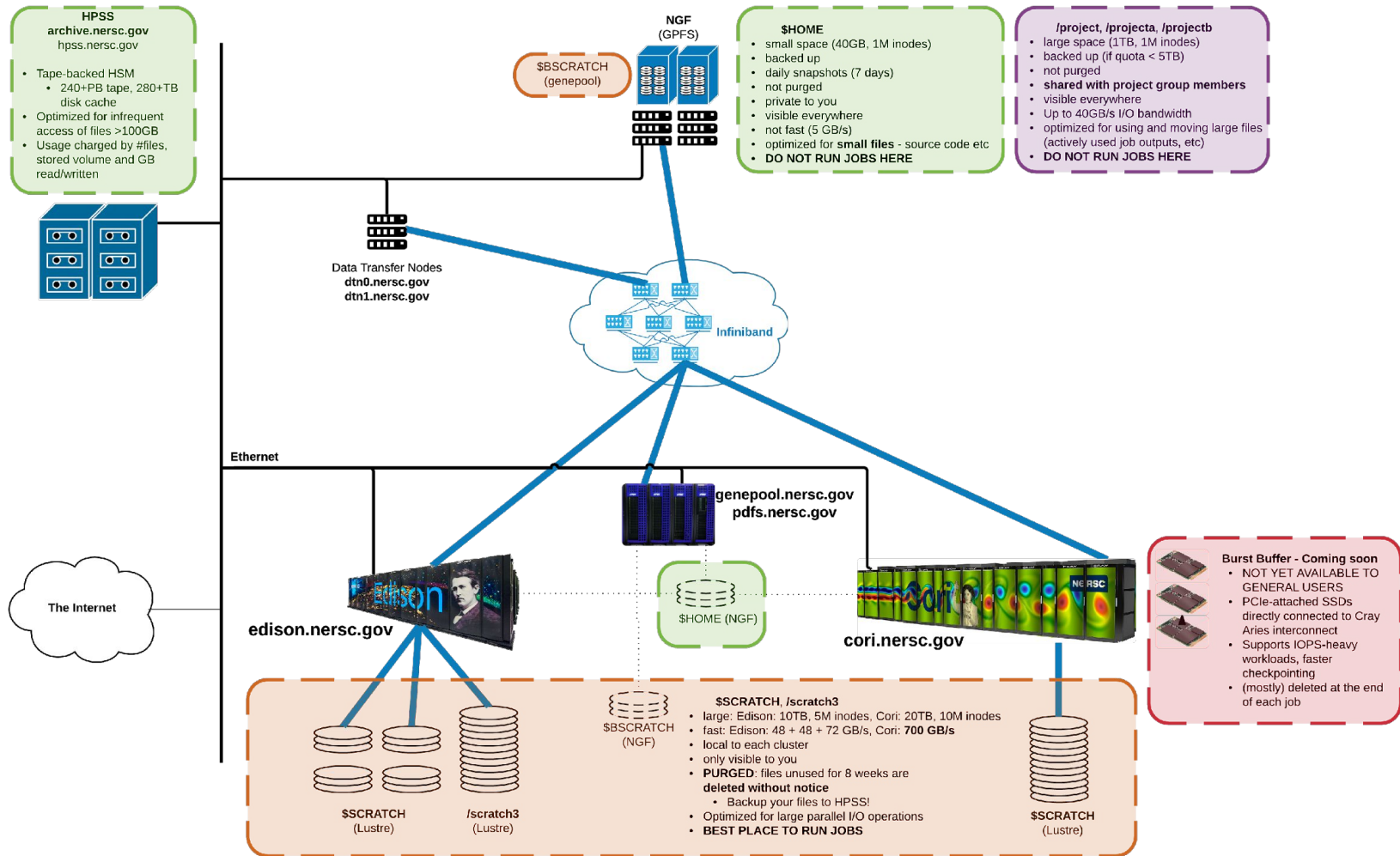
Moving Data Around



- **Don't do it!**
 - Ok, sometimes you need to
 - Don't forget \$HOME and /project are shared by all NERSC clusters
- **Data transfer nodes**
 - Fast network between all NERSC storage locations
 - Visible to internet
 - Dedicated to data transfer
 - Avoids adding load to Edison, Cori login nodes

NERSC File Systems Summary

NERSC



Moving Data Around



| Tool | What it does | Where/why to use it | Example |
|------------|---|--|---|
| cp | Local copy | Between NERSC filesystems | <pre>\$ cp \$SCRATCH/output.dat /project/projectdirs/m9999/</pre> |
| scp, rsync | Encrypted copy over network | Small amounts of data, collections of small files, over small distances. Use HPN version if available. | <pre>\$ scp my_code.f cori: \$ scp -R my_folder/ cori: \$ rsync -avr my_folder/ cori: \$ ssh -V OpenSSH_7.1p1-hpn14v5NMOD_3.17, OpenSSL 0.9.8j-fips 07 Jan 2009</pre> |
| bbcp | Fast parallel network copy. Requires client program | Larger files, longer distances | <pre>\$ bbcp -T "ssh -x -a -oFallbackToRsh=no %I -l %U %H /usr/common/usg/bin/bbcp" /local/path/file "user_name@dtm01.nersc.gov:/remote/path/"</pre> |

See <https://www.nersc.gov/users/storage-and-file-systems/transferring-data/>

Moving Data Around



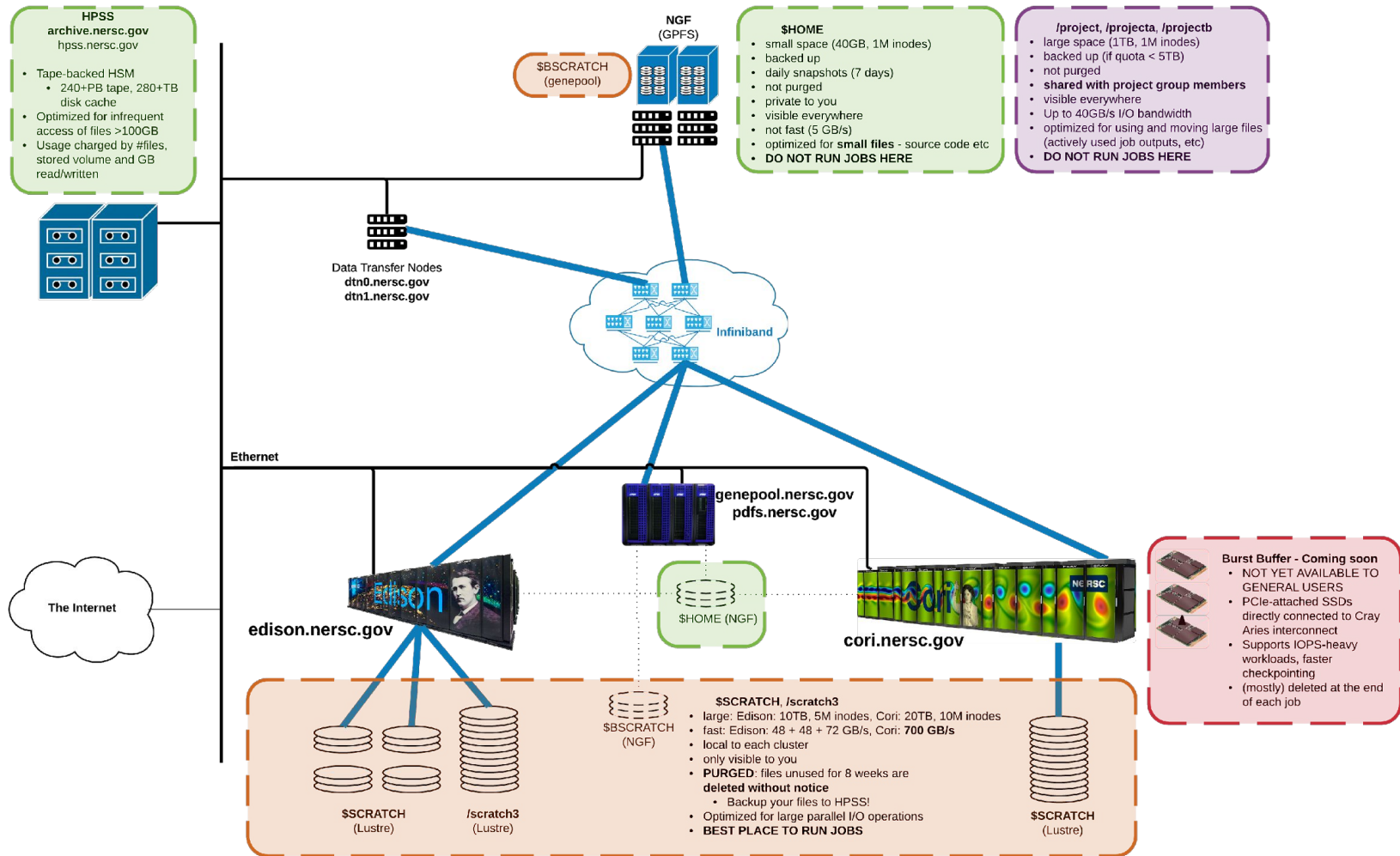
| Tool | What it does | Where/why to use it | Example |
|------------------|--|---|---|
| NERSC ftp upload | Temporary ftp account/server | Allow external collaborators to upload files for you to collect | See https://www.nersc.gov/users/storage-and-file-systems/transferring-data/nersc-ftp-upload-service/ |
| gridFTP | Fast network copy protocol, requires certificate | External, gridFTP-enabled sites (you need a grid credential) Note: garchiver.nersc.gov | <pre>\$ globus-url-copy file:///home/username/myresults.tar gsiftp://garchiver.nersc.gov/home/username/sleak/results-for-publication.tar</pre> |
| Globus Online | Fast data transfer service. Web or CLI | Fire-and-forget transfers (Especially between NERSC and other HPC centers) | See www.globusonline.org |

See <https://www.nersc.gov/users/storage-and-file-systems/transferring-data/>

- **Variety of storage types available to meet different needs**
 - Be aware of strengths and limitations of each, use each accordingly
- **BACK UP YOUR IMPORTANT FILES TO HPSS (archive)**
- **Many ways to move data to/from NERSC**
 - And most of them are better than 'scp'
- **If in doubt, ask for help**
 - www.nersc.gov -> "For Users"
 - ServiceNow (help.nersc.gov) or email (consult@nersc.gov)

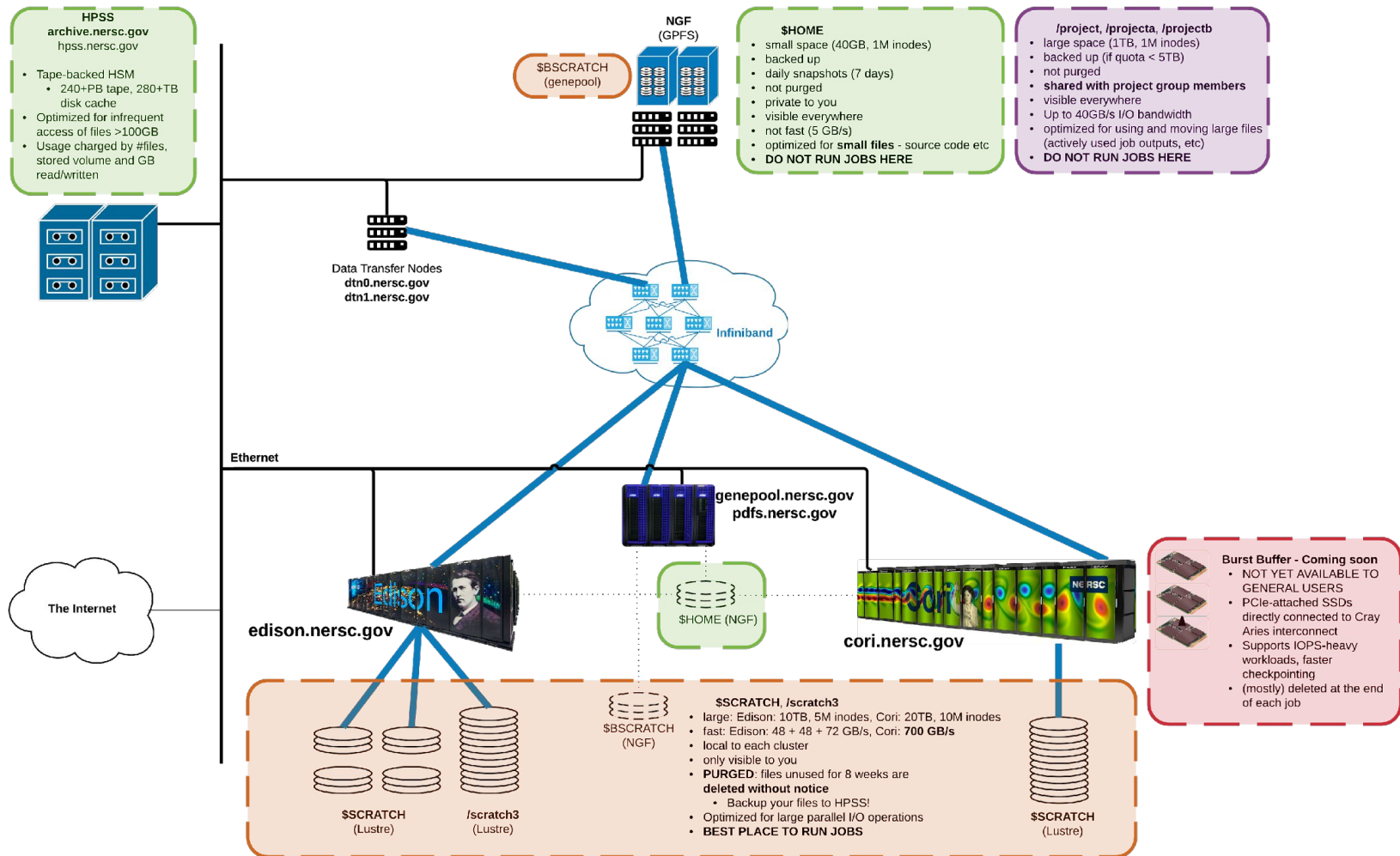
NERSC File Systems Summary

NERSC

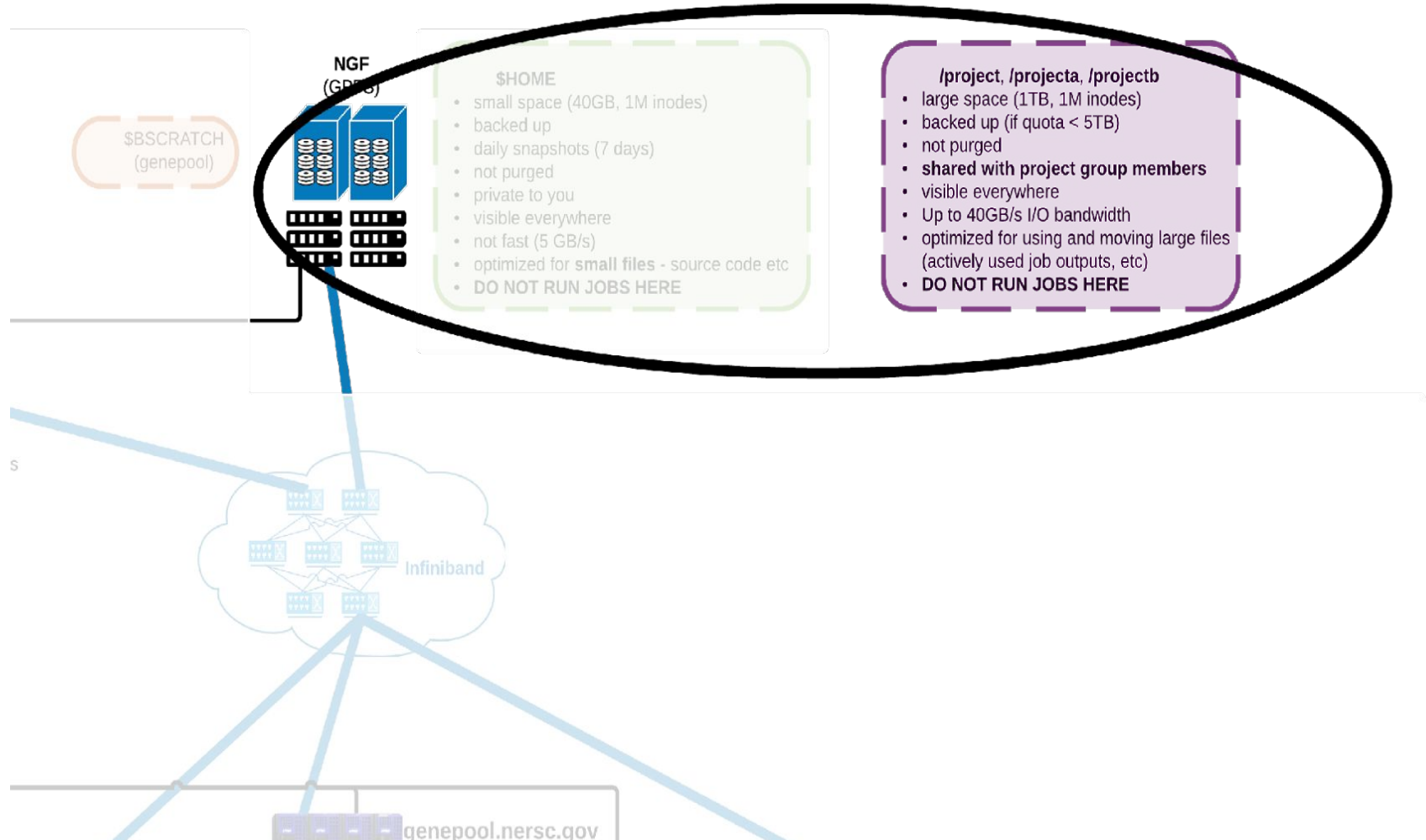


NERSC File Systems in a nutshell

NERSC

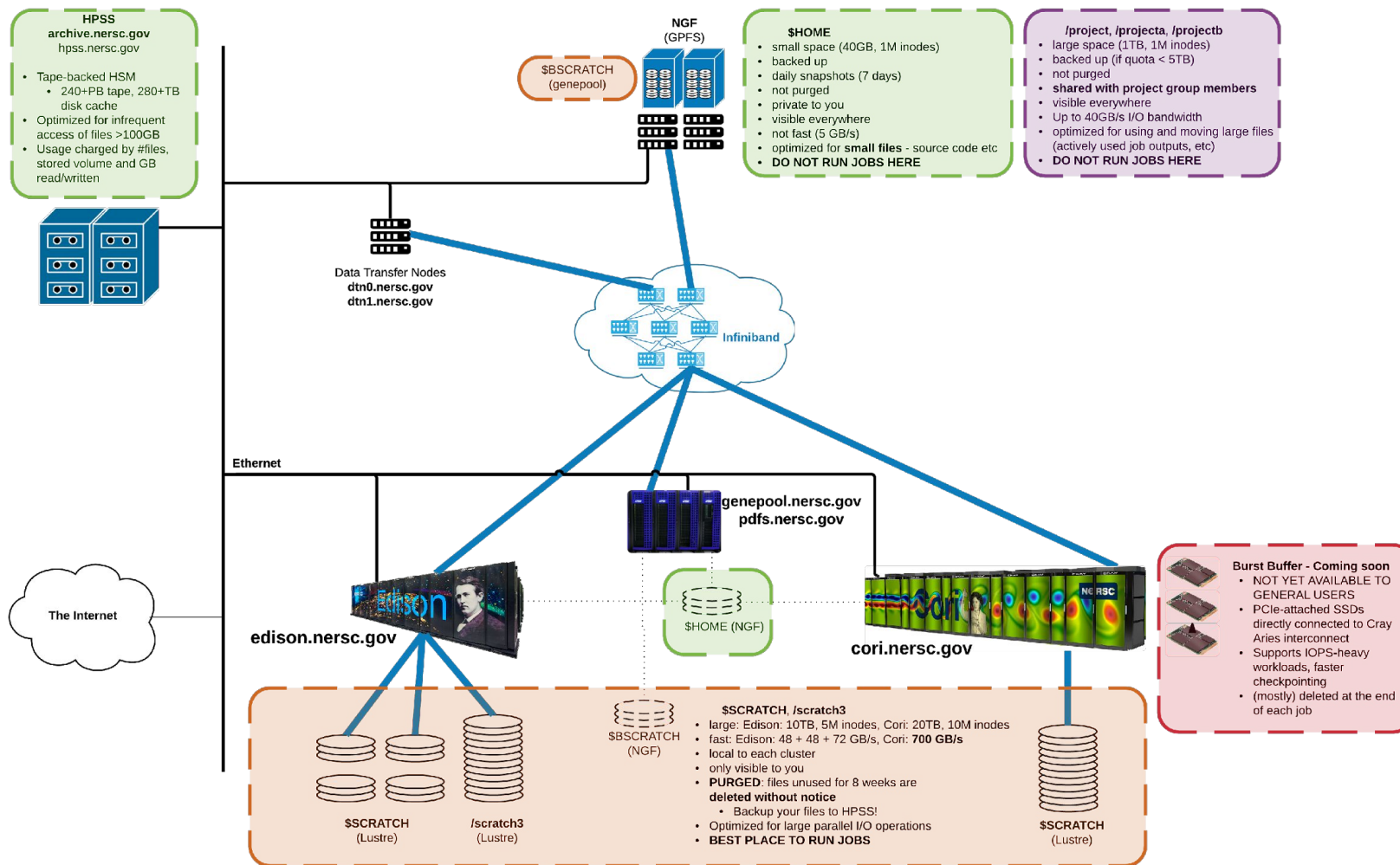


Project File Systems



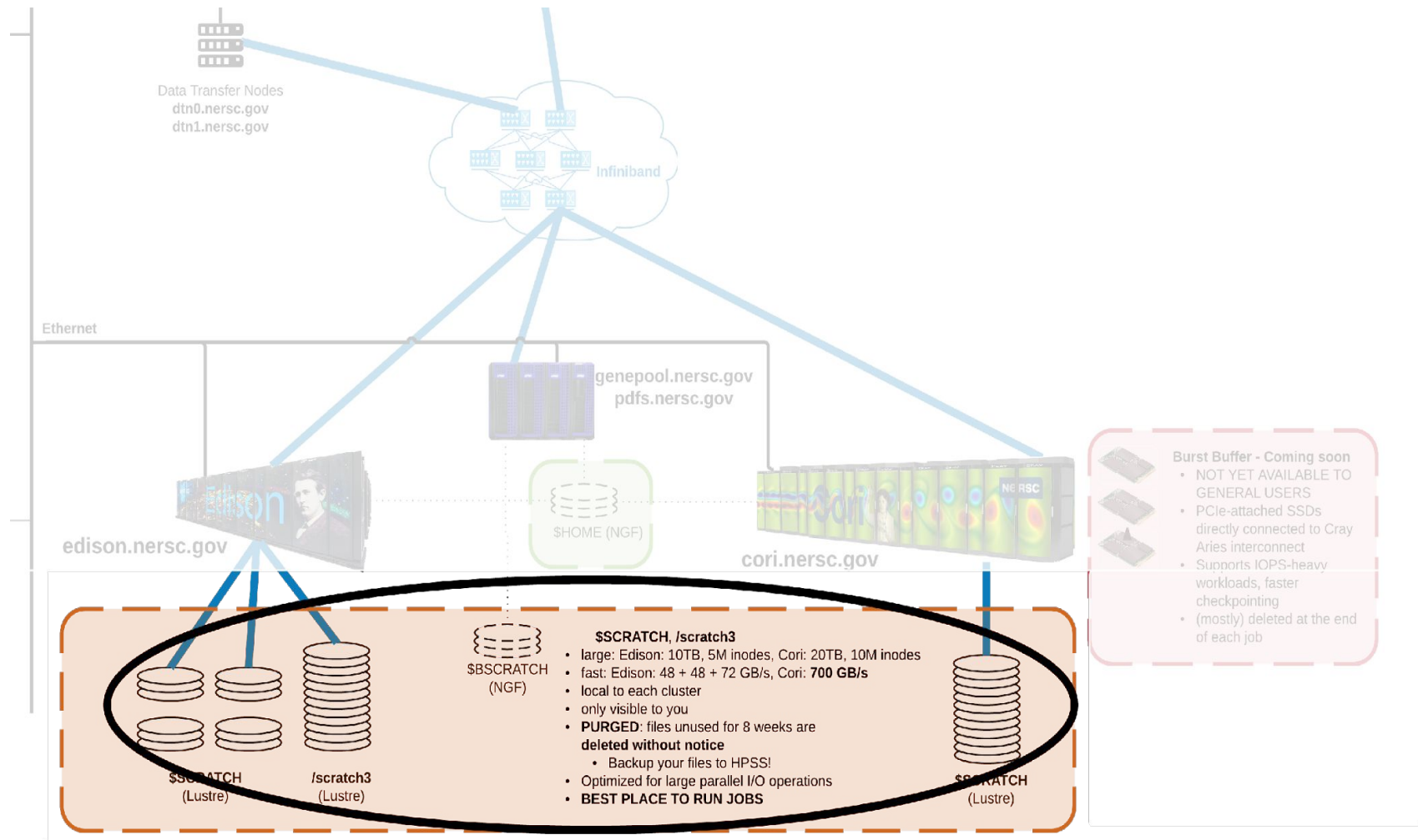
NERSC File Systems in a nutshell

NERSC



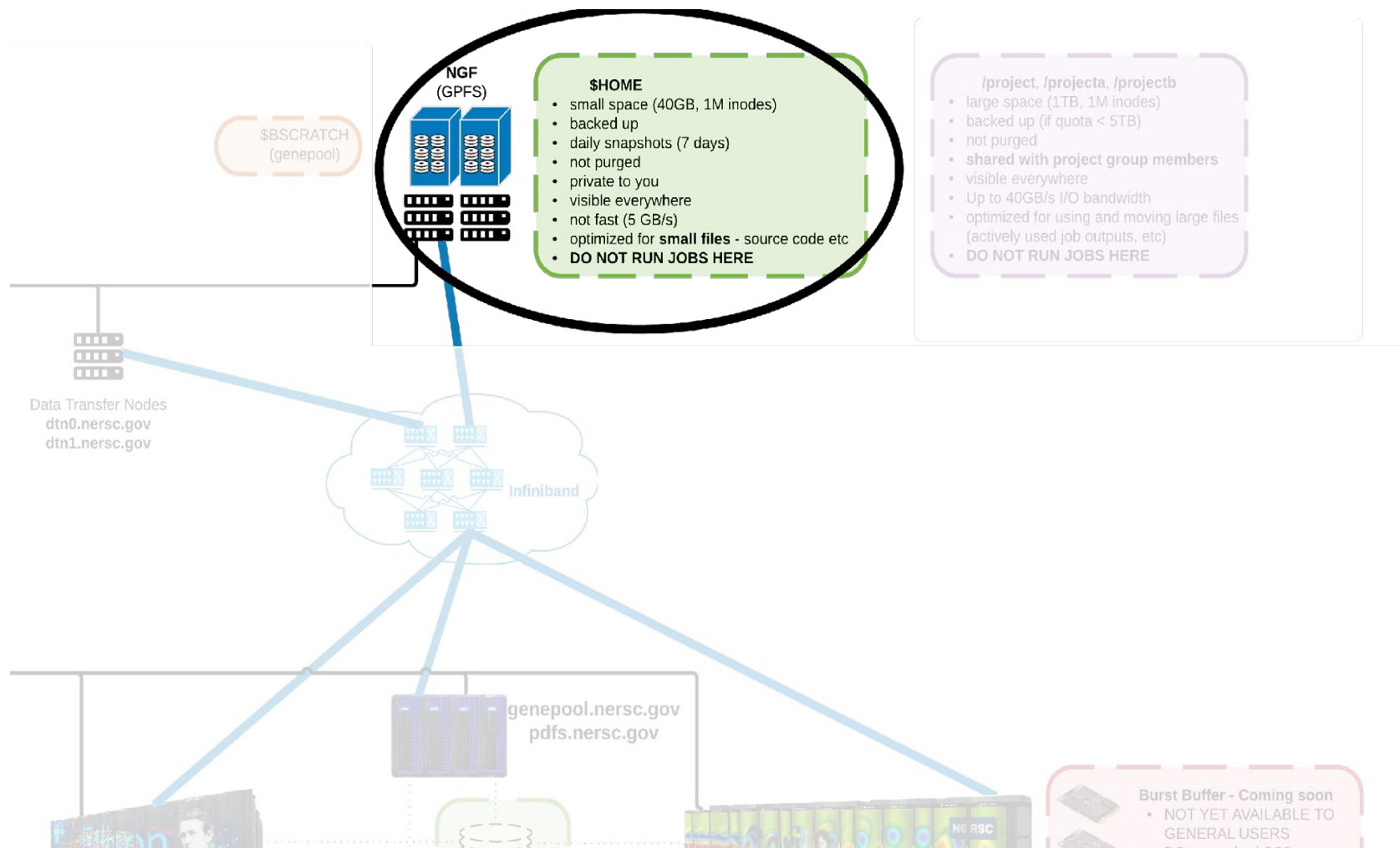
Local \$SCRATCH

NERSC



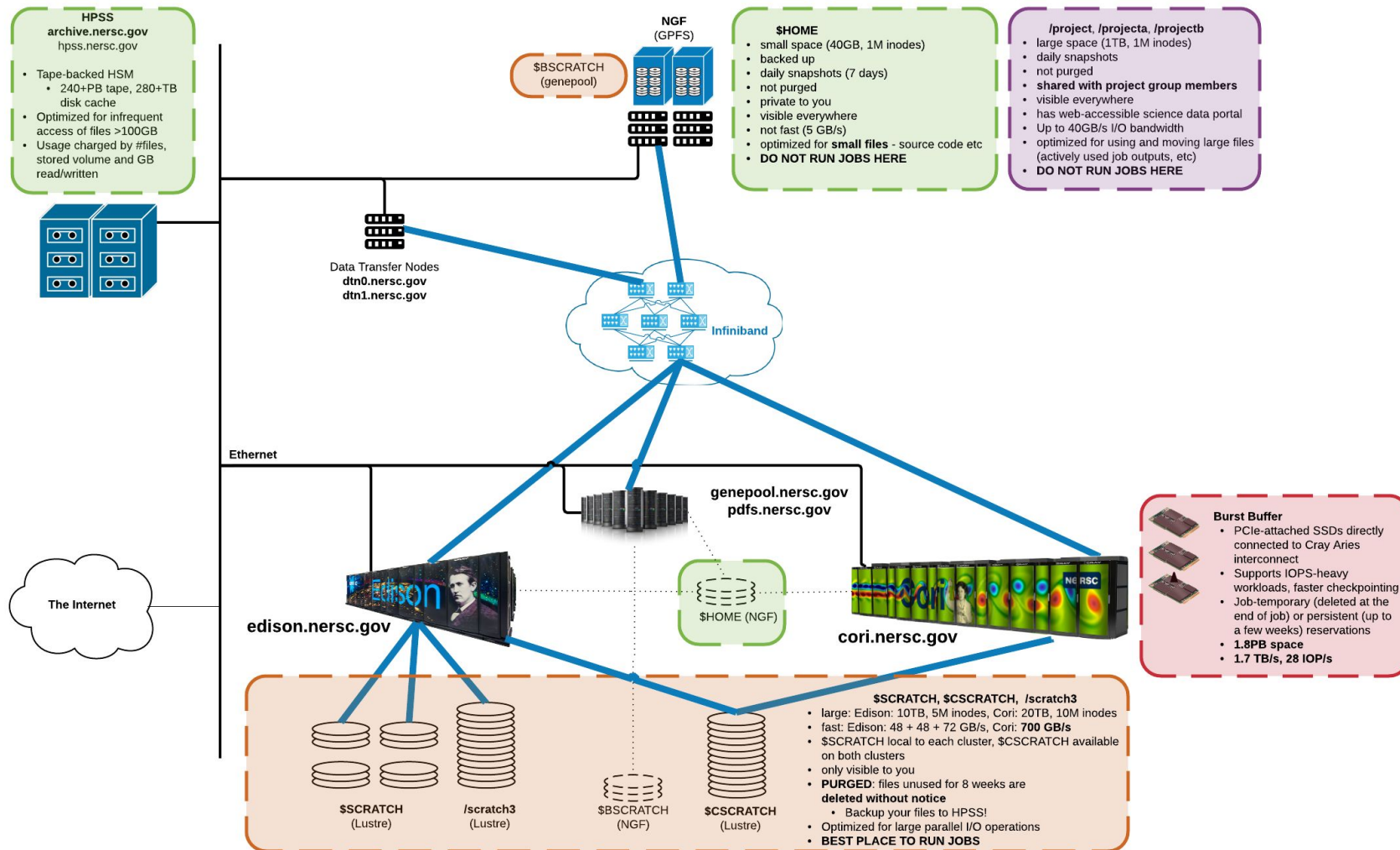
NERSC Global \$HOME

NERSC

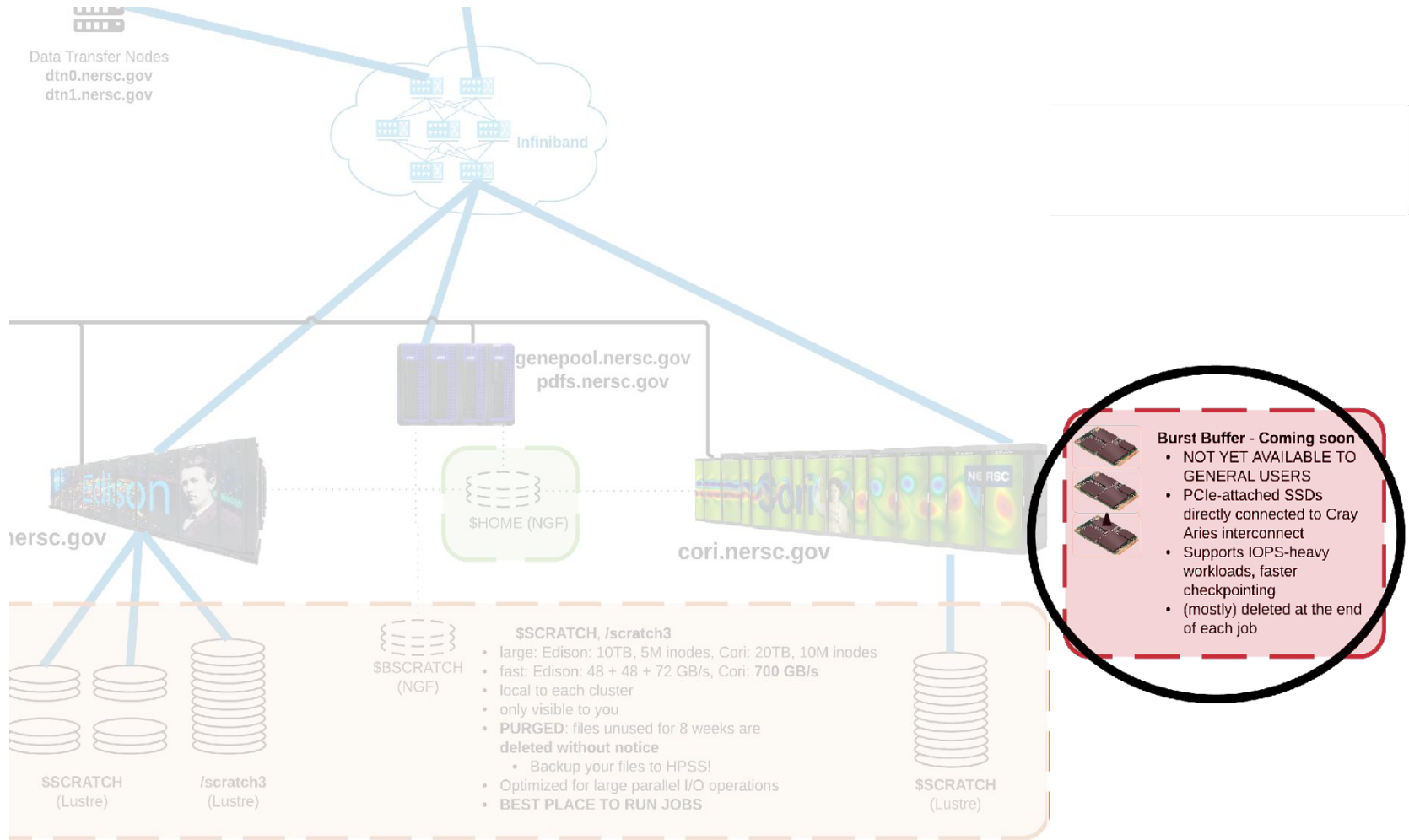


NERSC File Systems in a nutshell

NERSC



Burst Buffer



NERSC File Systems in a nutshell

NERSC

